



# MOLE-BLAST

A tool for clustering multiple sequences with their database neighbors

<https://www.ncbi.nlm.nih.gov/moleblast/moleblast.cgi>

National Center for Biotechnology Information • National Library of Medicine • National Institutes of Health • Department of Health and Human Services

## Overview

Large scale sequencing of unknown or environmental samples is a powerful method in biological research to help identify organisms present on the skin, in the gut, in soil sample, or in harsh environments such as toxic sites and natural ecosystems, to provide important insight on human health, bioremediation, global and local ecosystem processes. MOLE-BLAST is a tool designed cluster input sequences and place them in taxonomic context to help identify these sequences and the organisms they were isolated from. It works by taking multiple nucleotide sequences as input, searching them against a source database to find their best matches. The input query sequences along with their newly identified database matches are then subjected to multiple sequence alignment and the result is used to construct a phylogenetic tree. By placing input query sequences under the context of known sequences, MOLE-BLAST helps establish the taxonomic relationship to those known sequences and significantly increases the reliability of the identification process.

## Access

MOLE-BLAST is available from the “Specialized BLAST” section of the BLAST homepage ([blast.ncbi.nlm.nih.gov](http://blast.ncbi.nlm.nih.gov)). The search page is at [blast.ncbi.nlm.nih.gov/moleblast/moleblast.cgi](http://blast.ncbi.nlm.nih.gov/moleblast/moleblast.cgi).

The top input box (A) of this page accepts a group of query nucleotide sequences in FASTA format or a list of NCBI database identifiers. Clicking the “Align” button (B) starts the search with default settings with the neighboring sequences coming from NT nucleotide database. MOLE-BLAST also provides specialized databases including the 16S reference sequences from bacteria and archaea (C). Additional parameters for the search are available in the expandable “Advanced parameters” section (D). The “Database Search Parameters” section (E) controls the search sensitivity - megablast for faster searches, blastn for more sensitive searches, and number of matches to keep for the BLAST search step. It also allows database subset selection for more focused search. The “Query Clustering Parameters” section (F) sets thresholds for separating clusters into loci. These settings may separate sequences from different genes if they are present in the sample. The “Multiple Alignment Parameters” section (G) determines how many database matches found in the BLAST search step should be used in multiple alignment to construct the phylogenetic tree.

Pushing the “Align” button initiates the search and brings up the interim screen (H). **Leave this browser window open until the search is completed to avoid losing the assigned Request ID (RID, I) and the results.** Search results are save for 24 hours and can be retrieved using the assigned RID and the input box under the “Recent Results” tab.

Please, do not close the browser before the processing is completed	
Mole-BLAST Request ID	HC1DWV3J413
Status	Calculating
Submitted at	Fri Mar 27 21:36:55 2015
Current time	Fri Mar 27 21:38:45 2015
Time since submission	1 minute 50 seconds

## Example Searches

The first example is to establish the phylogeny of [a set sequences](#) from a clinical viral samples (A) using MOLE-BLAST. The search uses the NT database for its broad organism coverage, the Blastn program for increased sensitivity with “Max target sequences” set to 50 to increase the database reference matches saved. The “Show results in a new window” option is checked so input settings are available for subsequent adjustment and resubmission if necessary (B). Additional filtering of the database is done to ensure the quality of matches by excluding “Uncultured/environmental sample sequences,” selecting only sequences with “binomial” scientific names, and only database sequences from viruses are searched through the organism Entrez Query (C). The “Query Clustering Parameters” are left unchanged. The top 20 database sequences (D) out of the 50 saved in the BLAST steps are requested in the multiple sequence alignment and tree construction.

**MOLE-BLAST Neighbor Search Tool** Home Recent Results Help

Nucleotide

MOLE-BLAST searches for closest neighbors...

Enter Query Sequences

Enter nucleotide accessions, gis, or FASTA sequences (up to 300 input sequences with up to 5000 bases each)

Or, upload FASTA file Choose File No file chosen

Job Title

Choose Search Set

Database Nucleotide collection (nr/nt)

Align Show results in a new window

Advanced parameters

Database Search Parameters

Blast program Megablast Blastn

Max target sequences 50

Exclude Models (XM/XP) Uncultured/environmental sample sequences

Limit to Sequences from type material Sequences with a binomial name

Entrez Query viruses[orgn]

Query Clustering Parameters

Cluster queries Group query sequences by loci

Percent identity 40

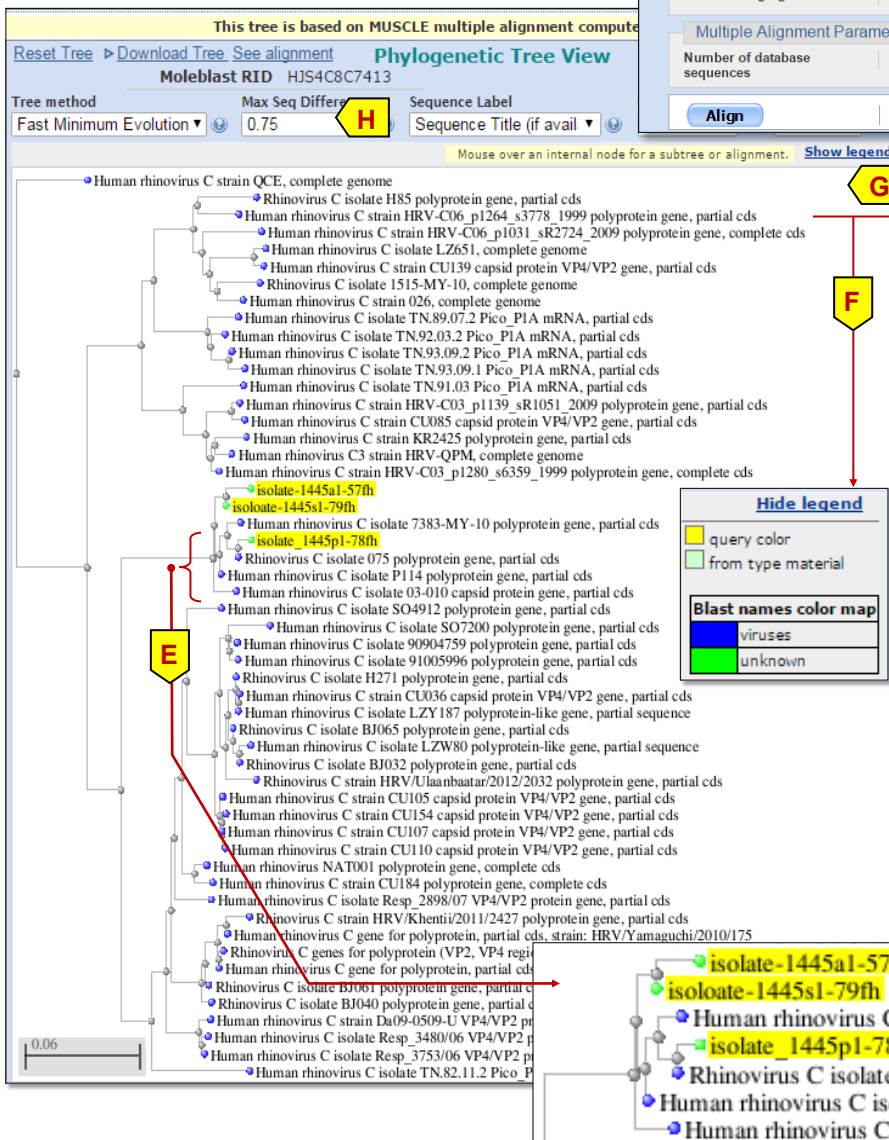
Percent sequence coverage 75

Cluster merging threshold 50%

Multiple Alignment Parameters

Number of database sequences 20

Align Show results in a new window



As shown to the left, MOLE-BLAST displays the result in TreeView format by default, with the query highlighted in yellow (E). The sequence color scheme indicated by the legend (F). The tree display can be adjusted by pull-down options at the top (G, from left to right): by selecting different methods, maximal sequence differences, different labels for sequences, collapsing or showing nodes, and showing different loci if the input queries are grouped into separate loci that cannot be merged. Links at the top of the page provide other functions (H, from left to right): clicking “Reset Tree” resets an adjusted tree display back to default; clicking the “Download Tree” displays link to tree files in different formats; and clicking the “See alignment” link displays the multiple sequence alignment, where the alignment can also be downloaded via a link at the top (not shown). The results indicate that the three sample are Human Rhinovirus C, closely related to previously reported isolates named as 7384-MY-10 and 075 (G).

## Example Searches (cont.)

The second example concerns identification of organisms present in wastewater samples, through 16S RNA sequences. [Input sequences](#) are uploaded through “Choose File” button and curated 16S RNA database is selected with other settings left at default (not shown).

In the resulting tree page (right), sequence titles are highlighted by background color to indicate if they are from type culture, with colored dot to indicate their major taxonomic group assignment (bacterial class), respectively. Color keys are listed in the Legend (A) shown to the right of the tree. The displayed group “Locus 1” contains eight of the eleven input queries (B). Results for other queries not displayed initially are available through the Locus pull down menu (C). Clicking the “See alignment” (D) link retrieves the underlying multiple alignment for the tree (E). The unit of the tree (F) represents expected number of changes per 100 bases. The tree is rooted at the middle of the longest edge.

**Phylogenetic Tree View**

This tree is based on MUSCLE multiple alignment computed for Mole-BLAST

Reset Tree ▶ Download Tree See alignment **D**

Moleblast RID: HP65KKN7413 Database nr: **C**

Tree method: Fast Minimum Evolution Max Seq Difference: 0.75 Sequence Label: Sequence Title (if avail) Show All Locus 1 Locus 2 Locus 3

Mouse over an internal node

**A** Legend

- query color
- from type material

**Blast names color map**

- a-proteobacteria
- unknown
- g-proteobacteria
- CFB group bacteria
- high GC Gram+

**B** Locus 1

**C** Locus 2

**D** See alignment

**E** Download alignment

Fasta plus gaps Clustal Phylip Nexus ASN.1

**F** 0.06

**G** Re-align

**H** Back to tree

**Descriptions**

Accession	Description	Links
NR_042275.1	Microcella alkaliphila strain AC4r 16S ribosomal RNA gene, complete sequence	
NR_042688.1	Leifsonia antarctica strain SPC-20 16S ribosomal RNA gene, complete sequence	
NR_042669.1	Leifsonia kafiensis strain KFC-22 16S ribosomal RNA gene, complete sequence	
NR_117545.1	Alpinimonas psychrophila strain Cr8-25 16S ribosomal RNA gene, partial sequence	
Ic Query_10008	APM129102:16-1530 Wastewater metagenome whole genome shotgun sequence	
NR_125643.1	Galbitalea soli strain KIS82-1 16S ribosomal RNA gene, partial sequence	
Ic Query_10004	APM153559:1-997 Wastewater metagenome whole genome shotgun sequence	
NR_074569.1	Gordonia sp. KTR9 strain KTR9 16S ribosomal RNA, complete sequence	
NR_043330.1	Gordonia rubripertincta strain DSM 43248 16S ribosomal RNA gene, partial sequence	
NR_104572.1	Gordonia rubripertincta strain N4 16S ribosomal RNA gene, complete sequence	
NR_025468.1	Gordonia westfalica strain Kb2 16S ribosomal RNA gene, partial sequence	

**Alignments**

Select All Re-align Mouse over the sequence identifier for sequence title

View Format: Plain Text

Accession	Sequence	Position
NR_042275	-----TGATCCTGGCTCAGGACGACGCTGGCGGCGTGCTTAACACATGCAA-GTCGAACGATGAA--C	61
NR_042688	-----AGAGTTTGATCCTGGCTCAGGACGACGCTGGCGGCGTGCTTAACACATGCAA-GTCGAACGATGAA--	66
NR_042669	-----TTTGAGTTTGATCCTGGCTCAGGACGACGCTGGCGGCGTGCTTAACACATGCAA-GTCGAACGATGAA--	68
NR_117545	-----AGAGTTTGATCCTGGCTCAGGACGACGCTGGCGGCGTGCTTAACACATGCAA-GTCGAACGATGAA--G	67
Query_10008	-----AGAGTTTGATCCTGGCTCAGGATGAACGCTGGCGGCGTGCTTAACACATGCAA-GTCGAACGATGAAAGC	69
NR_125643	-----ATGGCTCAGGATGAACGCTGGCGGCGTGCTTAACACATGCAA-GTCGAACGATGAAAGC	58
Query_10004	-----GAGTTTGATCCTGGCTCAGGACGACGCTGGCGGCGTGCTTAACACATGCAA-GTCGAACGAAAG--	65
NR_074569	-----TCAGGACGACGCTGGCGGCGTGCTTAACACATGCAA-GTCGAACGAAAG--	50
NR_043330	-----CCTGGCTCAGGACGACGCTGGCGGCGTGCTTAACACATGCAA-GTCGAACGAAAG--	56
NR_104572	-----GACGAACGCTGGCGGCGTGCTTAACACATGCAA-GTCGAACGAAAG--	46

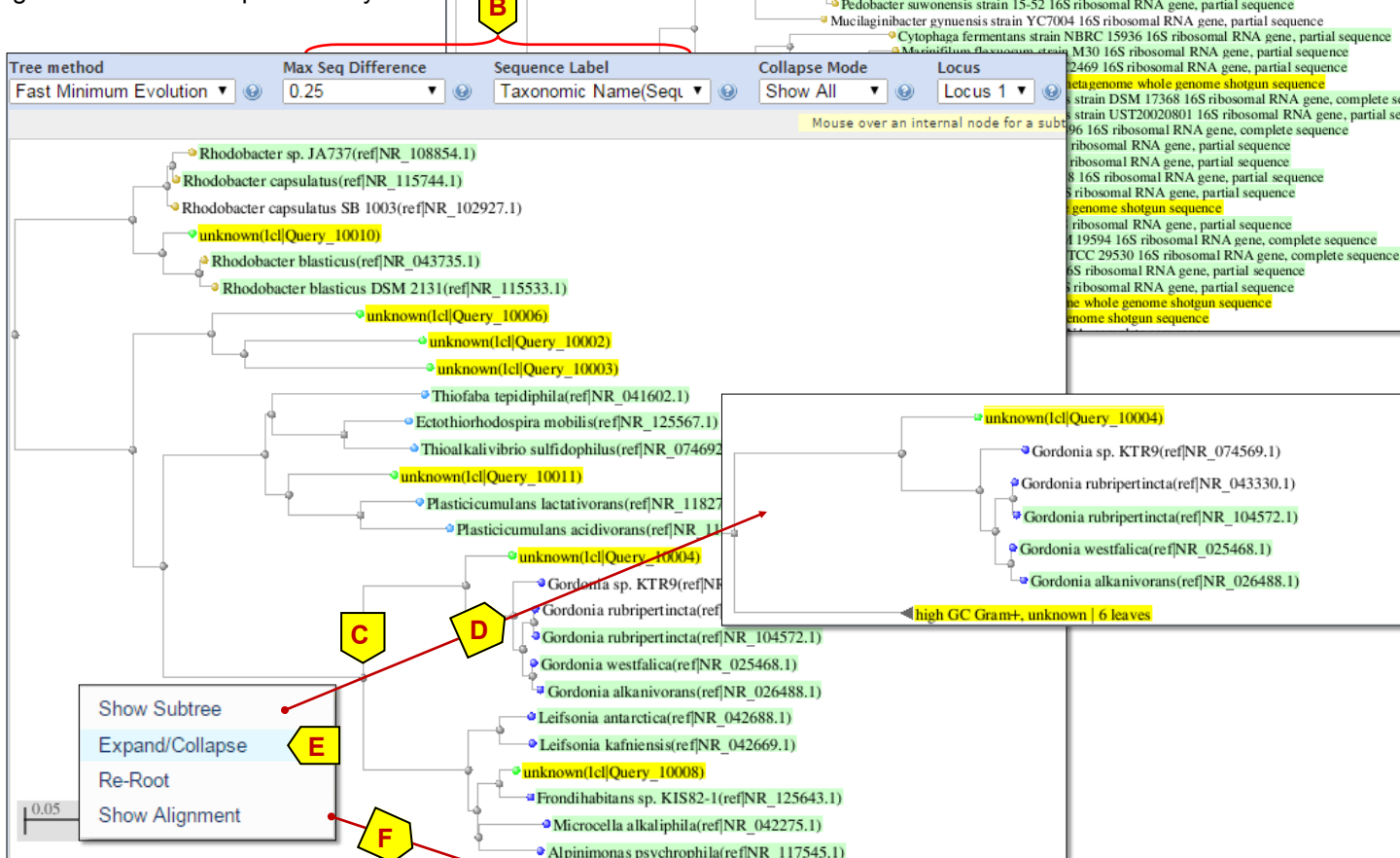
In the alignment view, the Descriptions section lists the sequence titles, and the Alignments section shows the details of aligned sequences (both shown in part to the left). The alignment can be redone using the “Re-align” button (G), after deselecting certain undesired entries using their checkboxes. The alignment file can be downloaded by clicking the “Download” link (H), which opens a section listing available formats.



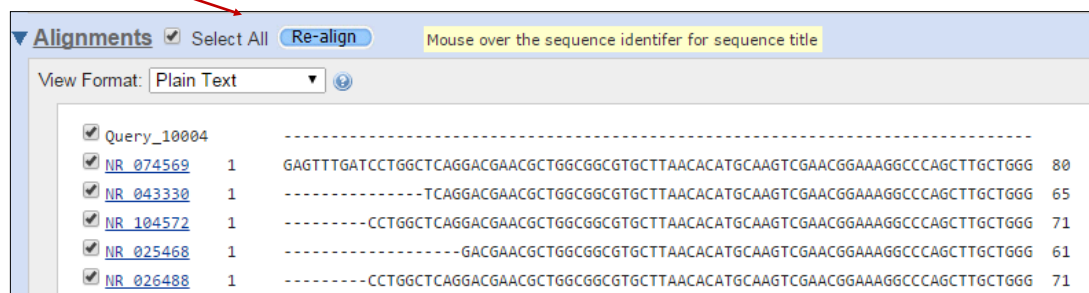
## Functions Provided by TreeView Display

A TreeView display provides two sets of functions to allow further customization of the result display.

One set is accessible through the pull-down menus listed at the top of the tree (A). In this case, reducing the “Max Seq Differences” to 0.25 (B) removes a query and a few database matches from the treeview display, and changing the “Sequence Label” to “Taxonomic Name (SeqId)” (C) shows the source organism name more prominently.



Another set of functions is accessible through the context menus within the tree display (C). There are activated upon hovering a non-leaf (internal) node of the tree. Selecting “Show Subtree” option (D) shows only that selected tree branch in a zoomed-in view. Selecting the “Expand/Collapse” option (E) collapses the branch under that node. Selecting the “Show Alignment” option (F) changes the display to show the alignment for that branch.



## Technical Assistance

Please send questions and feedback on MOLE-BLAST to the blast-help group at: [blast-help@ncbi.nlm.nih.gov](mailto:blast-help@ncbi.nlm.nih.gov)